



THE UNIVERSITY OF  
CHICAGO

# Stat 310 – Mihai Anitescu

## Lecture 8

---

## **8.4.1 OPTIMALITY CONDITIONS FOR EQUALITY CONSTRAINTS**

# IFT for optimality conditions in the equality-only case

- Problem:  $(NLP) \min f(x) \text{ subject to } c(x) = 0; c: \mathbb{R}^n \rightarrow \mathbb{R}^m$
- Assumptions:
  1.  $x^*$  is a solution
  2. LICQ:  $\nabla c(x)$  has full row rank.
- From LICQ:  $\exists x^* = \begin{pmatrix} \overbrace{x_D^*}^{n-m} \\ \overbrace{x_H^*}^m \end{pmatrix}; \nabla c_H(x^*) \in \mathbb{R}^{m \times m}; \nabla c_H(x^*) \text{ invertible.}$
- From IFT:
 
$$\exists \mathcal{N}(x^*), \Psi(x_D), \mathcal{N}(x_D^*) \text{ such that } x \in \mathcal{N}(x^*) \cap \Omega \Leftrightarrow x_H = \Psi(x_D)$$
- As a result  $x^*$  is a solution of NLP iff  $x_D^*$  solves unconstrained problem:  $\min_{x_D} f(x_D, \Psi(x_D))$

# Properties of Mapping

---

- From IFT:

$$c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})) = 0 \Rightarrow \nabla_{x_{\mathcal{D}}} c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})) + \nabla_{x_{\mathcal{H}}} c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})) \nabla_{x_{\mathcal{D}}} \Psi(x_{\mathcal{D}}) = 0$$

- Two important consequences

$$(1) \nabla_{x_{\mathcal{D}}} \Psi(x_{\mathcal{D}}) = - \left[ \nabla_{x_{\mathcal{H}}} c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})) \right]^{-1} \nabla_{x_{\mathcal{D}}} c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))$$

$$(2) Z = \begin{bmatrix} I_{n-m} \\ \nabla_{x_{\mathcal{D}}} \Psi(x_{\mathcal{D}}) \end{bmatrix} \Rightarrow \nabla c(x) Z = 0 \Rightarrow \text{Im}[Z] = \ker[\nabla c(x)]$$

# First-order optimality conditions

---

- Optimality of unconstrained optimization problem

$$\begin{aligned}\nabla_{x_D} f(x_D^*, \Psi(x_D^*)) &= 0 \Rightarrow \nabla_{x_D} f(x_D^*, \Psi(x_D^*)) + \nabla_{x_H} f(x_D^*, \Psi(x_D^*)) \nabla_{x_D} \Psi(x_D^*) = 0 \Rightarrow \\ \nabla_{x_D} f(x_D^*, \Psi(x_D^*)) - \underbrace{\nabla_{x_H} f(x_D^*, \Psi(x_D^*)) \left[ \nabla_{x_H} c(x_D, \Psi(x_D)) \right]^{-1}}_{\lambda^T} \nabla_{x_D} c(x_D, \Psi(x_D)) &= 0\end{aligned}$$

- The definition of the Lagrange Multiplier Result in the first-order (Lagrange, KKT) conditions:

$$\begin{aligned}\begin{bmatrix} \nabla_{x_D} f(x_D^*, \Psi(x_D^*)) & \nabla_{x_H} f(x_D^*, \Psi(x_D^*)) \end{bmatrix} - \lambda^T \begin{bmatrix} \nabla_{x_D} c(x_D, \Psi(x_D)) & \nabla_{x_H} c(x_D^*, \Psi(x_D^*)) \end{bmatrix} &= 0 \\ \nabla f(x^*) - \lambda^T \nabla c(x^*) &= 0\end{aligned}$$

# A more abstract and general proof

---

- Optimality of unconstrained optimization problem

$$D_{x_D} f(x_D^*, \Psi(x_D^*)) = 0 \Rightarrow \nabla_{x_D} f(x_D^*, \Psi(x_D^*)) + \nabla_{x_H} f(x_D^*, \Psi(x_D^*)) \nabla_{x_D} \Psi(x_D^*) = 0 \Rightarrow \nabla_x f(x^*) Z = 0$$

- Using  $\ker M \perp \text{Im } M^T$ ;  $\dim(\ker M) + \dim(\text{Im } M^T) = \text{nr cols } M$
- We obtain:  $\nabla_x f(x^*) Z = 0 \Rightarrow \nabla_x f(x^*)^T \in \ker(Z^T) = \text{Im}[\nabla c(x^*)^T]$
- We thus obtain the optimality conditions:

$$\exists \lambda \in \mathbb{R}^m \text{ s.t. } \nabla_x f(x^*)^T = \nabla_x c(x^*)^T \lambda \Rightarrow \nabla_x f(x^*) - \lambda^T \nabla_x c(x^*) = 0$$

# The Lagrangian

---

- Definition  $\mathcal{L}(x, \lambda) = f(x) - \lambda^T c(x)$
- Its gradient  $\nabla \mathcal{L}(x, \lambda) = \begin{bmatrix} \nabla f(x) - \lambda^T \nabla c(x), & c(x)^T \end{bmatrix}$
- Its Hessian  $\nabla^2 \mathcal{L}(x, \lambda) = \begin{bmatrix} \nabla_{xx}^2 \mathcal{L}(x, \lambda) & \nabla c(x)^T \\ \nabla c(x) & 0 \end{bmatrix}$
- Where  $\nabla_{xx}^2 \mathcal{L}(x, \lambda) = \nabla_{xx}^2 f(x, \lambda) - \sum_{i=1}^m \lambda_i \nabla_{xx}^2 c_i(x, \lambda)$
- Optimality conditions:  $\boxed{\nabla \mathcal{L}(x, \lambda) = 0}$

# Second-order conditions

---

- First, note that:  $Z^T \nabla_{xx}^2 L(x_D, \Psi(x_D)) Z = D_{x_D x_D}^2 f(x_D, \Psi(x_D)) \succcurlyeq 0$
- Sketch of proof: total derivatives in  $x_D$  :

$$\begin{aligned} D_{x_D} f(x_D, \Psi(x_D)) &= \nabla_{x_D} f(x_D, \Psi(x_D)) - \lambda(x_D, \Psi(x_D))^T \nabla_{x_D} c(x_D^*, \Psi(x_D^*)) = \\ &\nabla_{x_D} \mathcal{L}((x_D, \Psi(x_D)), \lambda(x_D, \Psi(x_D))); \\ \nabla_{x_{\mathcal{H}}} f(x_D^*, \Psi(x_D^*)) &= \lambda(x_D, \Psi(x_D))^T \nabla_{x_{\mathcal{H}}} c(x_D^*, \Psi(x_D^*)) \end{aligned}$$

- Second derivatives:

$$\begin{aligned} D_{x_D x_D} f(x_D, \Psi(x_D)) &= \nabla_{x_D} f(x_D, \Psi(x_D)) - \lambda(x_D, \Psi(x_D))^T \nabla_{x_D} c(x_D, \Psi(x_D)) = \\ \nabla_{x_D x_D} \mathcal{L}((x_D, \Psi(x_D)), \lambda(x_D, \Psi(x_D))) &+ \nabla_{x_D} \Psi(x_D)^T \nabla_{x_{\mathcal{H}} x_D} \mathcal{L}((x_D, \Psi(x_D)), \lambda(x_D, \Psi(x_D))) \\ - D_D \left( \lambda(x_D, \Psi(x_D))^T \right) \nabla_{x_D} c(x_D, \Psi(x_D)) \end{aligned}$$



# Computing Second-Order Derivatives

---

- Expressing the second derivatives of Lagrangian

$$\begin{aligned} \nabla_{x_{\mathcal{H}}} f(x_{\mathcal{D}}^*, \Psi(x_{\mathcal{D}}^*)) &= \lambda(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))^T \nabla_{x_{\mathcal{H}}} c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})) \Rightarrow \\ D_{x_{\mathcal{D}}} \left[ \lambda(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))^T \right] \nabla_{x_{\mathcal{H}}} c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})) &= D_{x_{\mathcal{D}}} \left[ \nabla_{x_{\mathcal{H}}} f(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})) - \underbrace{\lambda(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))^T \nabla_{x_{\mathcal{H}}} c(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))}_{inactive} \right] = \\ D_{x_{\mathcal{D}}} \nabla_{x_{\mathcal{H}}} \mathcal{L} \left( (x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})), \underbrace{\lambda(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))^T}_{inactive} \right) &= \nabla_{x_{\mathcal{D}}} \nabla_{x_{\mathcal{H}}} \mathcal{L} \left( (x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})), \lambda(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))^T \right) + \\ \nabla_{x_{\mathcal{D}}} \Psi(x_{\mathcal{D}})^T \nabla_{x_{\mathcal{H}}} \nabla_{x_{\mathcal{H}}} \mathcal{L} \left( (x_{\mathcal{D}}, \Psi(x_{\mathcal{D}})), \lambda(x_{\mathcal{D}}, \Psi(x_{\mathcal{D}}))^T \right) \end{aligned}$$

- Solve for total derivative of multiplier and replace conclusion follows.

# Summary: Necessary Optimality Conditions

---

- Summary:

$$\boxed{\nabla \mathcal{L}(x^*, \lambda^*) = 0; \quad Z^T \nabla_{xx}^2 L(x_D^*, \Psi(x_D^*)) Z \succcurlyeq 0}$$

- Rephrase first order:

$$\nabla_x \mathcal{L}(x^*, \lambda^*) = 0$$

- Rephrase second order necessary conditions.

$$\nabla_x c(x^*) w = 0 \Rightarrow w^T \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) w \geq 0$$

# Sufficient Optimality Conditions

---

- The point is a local minimum if LICQ and the following holds:

$$(1) \nabla_x \mathcal{L}(x^*, \lambda^*) = 0; (2) \nabla_x c(x^*) w = 0 \Rightarrow \exists \sigma > 0 \ w^T \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) w \geq \sigma \|w\|^2$$

- Proof: By IFT, there is a change of variables such that

$$u \in \mathcal{N}(0) \subset \mathbb{R}^{n-n_c} u \leftrightarrow x(u); \tilde{x} \in \mathcal{N}(x^*), c(\tilde{x}) = 0 \Leftrightarrow \exists \tilde{u} \in \mathcal{N}(0); \tilde{x} = x(\tilde{u})$$

$$\nabla_x c(x^*) \nabla_u x(\tilde{u}) \Big|_{\tilde{u}=0} = 0; \quad Z = \nabla_u x(\tilde{u})$$

- The original problem can be phrased as

$$\min_u f(x(u))$$

# Sufficient Optimality Conditions

---

- We can now piggy back on theory of unconstrained optimization, noting that.

$$\nabla_u f(x(u))\big|_{u=0} = \nabla_x \mathcal{L}(x^*, \lambda^*) = 0;$$

$$\nabla_{uu}^2 f(x(u))\big|_{u=0} = Z^T \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) Z \succ 0; Z = \nabla_u x(u)$$

- Then from theory of unconstrained optimization we have a local isolated minimum at 0 and thus the original problem at  $x^*$ . (following the local isomorphism above)

# Another Essential Consequence

---

- If LICQ+ second-order conditions hold at the solution  $x^*$ , then the following matrix must be nonsingular (**EXPAND**).

$$\begin{bmatrix} \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) & \nabla_x c(x^*) \\ \nabla_x^T c(x^*) & 0 \end{bmatrix}$$

- The system of nonlinear equations has an invertible Jacobian,

$$\begin{bmatrix} \nabla_x \mathcal{L}(x^*, \lambda^*) \\ c(x^*) \end{bmatrix} = 0$$

---

## 8.4.2 FIRST-ORDER OPTIMALITY CONDITIONS FOR MIXED EQ AND INEQ CONSTRAINTS

# The Lagrangian

---

- Even in the general case, it has the same expression

$$\mathcal{L}(x) = f(x) - \sum_{i \in \mathcal{C} \cup \mathcal{A}} \lambda_i c_i(x)$$

# First-Order Optimality Condition

---

## Theorem

*Suppose that  $x^*$  is a local solution of (12.1), that the functions  $f$  and  $c_i$  in (12.1) are continuously differentiable, and that the LICQ holds at  $x^*$ . Then there is a Lagrange multiplier vector  $\lambda^*$ , with components  $\lambda_i^*$ ,  $i \in \mathcal{E} \cup \mathcal{I}$ , such that the following conditions are satisfied at  $(x^*, \lambda^*)$*

$$\nabla_x \mathcal{L}(x^*, \lambda^*) = 0, \quad (12.34a)$$

$$c_i(x^*) = 0, \quad \text{for all } i \in \mathcal{E}, \quad (12.34b)$$

$$c_i(x^*) \geq 0, \quad \text{for all } i \in \mathcal{I}, \quad (12.34c)$$

$$\lambda_i^* \geq 0, \quad \text{for all } i \in \mathcal{I}, \quad (12.34d)$$

$$\lambda_i^* c_i(x^*) = 0, \quad \text{for all } i \in \mathcal{E} \cup \mathcal{I}. \quad (12.34e)$$

Equivalent Form:

$$\nabla f(x^*) - \lambda_{\mathcal{A}(x^*)}^T \nabla c_{\mathcal{A}(x^*)}(x^*) = 0 \Rightarrow \text{Multipliers are unique !!}$$



# Sketch of the Proof

---

- If  $x^*$  is a solution of the original problem, it is also a solution of the problem.

$$\min f(x) \text{ subject to } c_{\mathcal{A}(x^*)}(x) = 0$$

- From the optimality conditions of the problem with equality constraints, we must have (since LICQ holds)

$$\exists \{\lambda_i\}_{i \in \mathcal{A}(x^*)} \text{ such that } \nabla f(x^*) - \sum_{i \in \mathcal{A}(x^*)} \lambda_i \nabla c_i(x^*) = 0$$

- But I cannot yet tell by this argument  $\lambda_i \geq 0$

# Sketch of the Proof: The sign of the multiplier

---

- Assume now one multiplier has the “wrong” sign. That is
$$j \in \mathcal{A}(x^*) \cap \mathcal{I}, \quad \lambda_j < 0$$
- Since LICQ holds, we can construct a feasible path that “takes off” from that constraint (inactive constraints do not matter locally)

- $c_{\mathcal{A}(x^*)}(\tilde{x}(t)) = te_j \Rightarrow \tilde{x}(t) \in \Omega$  Define  $b = \frac{d}{dt} \tilde{x}(t)_{t=0} \Rightarrow \nabla c_{\mathcal{A}(x)} b = e_j$ 
$$\frac{d}{dt} f(\tilde{x}(t))_{t=0} = \nabla f(x^*)^T b = \lambda_{c_{\mathcal{A}(x)}}^T \nabla c_{\mathcal{A}(x)} b = \lambda_j < 0 \Rightarrow$$
$$\exists t_1 > 0, \quad f(\tilde{x}(t_1)) < f(\tilde{x}(0)) = f(x^*), \quad \text{CONTRADICTION!!}$$

# Strict Complementarity

---

- It is a notion that makes the problem look “almost” like an equality.

**Definition 12.5** (Strict Complementarity).

*Given a local solution  $x^*$  of (12.1) and a vector  $\lambda^*$  satisfying (12.34), we say that the strict complementarity condition holds if exactly one of  $\lambda_i^*$  and  $c_i(x^*)$  is zero for each index  $i \in \mathcal{I}$ . In other words, we have that  $\lambda_i^* > 0$  for each  $i \in \mathcal{I} \cap \mathcal{A}(x^*)$ .*

---

## **8.5 SECOND-ORDER CONDITIONS**

# Critical Cone

---

- The subset of the tangent space, where the objective function does not vary to first-order.
- The book definition.

$$\mathcal{C}(x^*, \lambda^*) = \{w \in \mathcal{F}(x^*) \mid \nabla c_i(x^*)^T w = 0, \text{ all } i \in \mathcal{A}(x^*) \cap \mathcal{I} \text{ with } \lambda_i^* > 0\}.$$

- An even simpler equivalent definition.

$$\mathcal{C}(x^*, \lambda^*) = \left\{ w \in T_{\Omega}(x^*) \mid \nabla f(x^*)^T w = 0 \right\}$$

# Rephrasing of the Critical Cone

---

- By investigating the definition

$$w \in \mathcal{C}(x^*, \lambda^*) \Leftrightarrow \begin{cases} \nabla c_i(x^*)^T w = 0 & i \in \mathcal{E} \\ \nabla c_i(x^*)^T w = 0 & i \in \mathcal{A}(x^*) \cap \mathcal{I} \quad \lambda_i^* > 0 \\ \nabla c_i(x^*)^T w \geq 0 & i \in \mathcal{A}(x^*) \cap \mathcal{I} \quad \lambda_i^* = 0 \end{cases}$$

- In the case where strict complementarity holds, the cone has a MUCH simpler expression.

$$w \in \mathcal{C}(x^*, \lambda^*) \Leftrightarrow \nabla c_i(x^*)^T w = 0 \quad \forall i \in \mathcal{A}(x^*)$$

# Statement of the Second-Order Conditions

---

**Theorem 12.5** (Second-Order Necessary Conditions).

*Suppose that  $x^*$  is a local solution of (12.1) and that the LICQ condition is satisfied. Let  $\lambda^*$  be the Lagrange multiplier vector for which the KKT conditions (12.34) are satisfied. Then*

$$w^T \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) w \geq 0, \quad \text{for all } w \in \mathcal{C}(x^*, \lambda^*). \quad (12.57)$$

- How to prove this? In the case of Strict Complementarity the critical cone is the same as the problem constrained with equalities on active index.
- Result follows from equality-only case.

# Statement of second-order sufficient conditions

---

**Theorem 12.6** (Second-Order Sufficient Conditions).

*Suppose that for some feasible point  $x^* \in \mathbb{R}^n$  there is a Lagrange multiplier vector  $\lambda^*$  such that the KKT conditions (12.34) are satisfied. Suppose also that*

$$w^T \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) w > 0, \quad \text{for all } w \in \mathcal{C}(x^*, \lambda^*), w \neq 0. \quad (12.65)$$

*Then  $x^*$  is a strict local solution for (12.1).*

- How do we prove this? In the case of strict complementarity again from reduction to the equality case.

$$x^* = \arg \min_x f(x) \text{ subject to } c_A(x) = 0$$



# How to derive those conditions in the other case?

---

- Use the slacks to reduce the problem to one with equality constraints.

$$\begin{array}{ll}\min_{x \in \mathbb{R}^n, z \in \mathbb{R}^{n_I}}, & f(x) \\ \text{s.t.} & c_E(x) = 0 \\ & [c_I(x)]_j - z_j^2 = 0 \quad j = 1, 2, \dots, n_I\end{array}$$

- Then, apply the conditions for equality constraints.
- I will assign it as homework.

# Summary: Why should I care about Lagrange Multipliers?

---

- Because it makes the optimization problem in principle equivalent to a nonlinear equation.

$$\begin{bmatrix} \nabla_x \mathcal{L}(x^*, \lambda^*) \\ c_{\mathcal{A}}(x^*) \end{bmatrix} = 0; \quad \det \begin{bmatrix} \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) & \nabla_x c_{\mathcal{A}}(x^*) \\ \nabla_x^T c_{\mathcal{A}}(x^*) & 0 \end{bmatrix} \neq 0$$

- I can use concepts from nonlinear equations such as Newton's for the algorithmics.



THE UNIVERSITY OF  
CHICAGO

# Section 9

## Fundamentals of Algorithms for Constrained Optimization

---

## **9.1 TYPES OF CONSTRAINED OPTIMIZATION ALGORITHMS**

# Types of Optimization Algorithms

---

- All of the algorithms solve iteratively a simpler problem.
  - Penalty and Augmented Lagrangian Methods.
  - Sequential Quadratic Programming.
  - Interior-point Methods.
- The approach follows the usual divide-and-conquer approach:
  - Constrained Optimization-
  - Unconstrained Optimization
  - Nonlinear Equations
  - Linear Equations

# Quadratic Programming Problems

---

- Algorithms for such problems are interested to explore because
  - 1. Their structure can be efficiently exploited.
  - 2. They form the basis for other algorithms, such as augmented Lagrangian and Sequential quadratic programming problems.

$$\begin{array}{ll} \min_x & q(x) = \frac{1}{2}x^T Gx + x^T c \\ \text{subject to} & a_i^T x = b_i, \quad i \in \mathcal{E}, \\ & a_i^T x \geq b_i, \quad i \in \mathcal{I}, \end{array}$$

# Penalty Methods

---

- Idea: Replace the constraints by a penalty term.
- Inexact penalties: parameter driven to infinity to recover solution. Example:

$$x^* = \arg \min f(x) \text{ subject to } c(x) = 0 \Leftrightarrow$$

$$x^\mu = \arg \min f(x) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} c_i^2(x); \quad x^* = \lim_{\mu \rightarrow \infty} x^\mu = x^*$$

Solve with unconstrained optimization

- Exact but nonsmooth penalty – the penalty parameter can stay finite.

$$x^* = \arg \min f(x) \text{ subject to } c(x) = 0 \Leftrightarrow x^* = \arg \min f(x) + \mu \sum_{i \in \mathcal{E}} |c_i(x)|; \quad \mu \geq \mu_0$$

# Augmented Lagrangian Methods

---

- Mix the Lagrangian point of view with a penalty point of view.

$$x^* = \arg \min f(x) \text{ subject to } c(x) = 0 \Leftrightarrow$$

$$x^{\mu, \lambda} = \arg \min f(x) - \sum_{i \in \mathcal{E}} \lambda_i c_i(x) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} c_i^2(x) \Rightarrow$$

$$x^* = \lim_{\lambda \rightarrow \lambda^*} x^{\mu, \lambda} \text{ for some } \mu \geq \mu_0 > 0$$



# Sequential Quadratic Programming

---

## Algorithms

- Solve successively Quadratic Programs.

$$\begin{aligned} \min_p \quad & \frac{1}{2} p^T B_k p + \nabla f(x_k) \\ \text{subject to} \quad & \nabla c_i(x_k) d + c_i(x_k) = 0 \quad i \in \mathcal{E} \\ & \nabla c_i(x_k) d + c_i(x_k) \geq 0 \quad i \in \mathcal{I} \end{aligned}$$

- It is the analogous of Newton's method for the case of constraints if  $B_k = \nabla_{xx}^2 \mathcal{L}(x_k, \lambda_k)$
- But how do you solve the subproblem? It is possible with extensions of simplex which I do not cover.
- An option is BFGS which makes it convex.

# Interior Point Methods

---

- Reduce the inequality constraints with a barrier

$$\begin{aligned} \min_{x,s} \quad & f(x) - \mu \sum_{i=1}^m \log s_i \\ \text{subject to} \quad & c_i(x) = 0 \quad i \in \mathcal{E} \\ & c_i(x) - s_i = 0 \quad i \in \mathcal{I} \end{aligned}$$

- An alternative, is use a penalty as well:

$$\min_x f(x) - \mu \sum_{i \in \mathcal{I}} \log s_i + \frac{1}{2\mu} \sum_{i \in \mathcal{I}} (c_i(x) - s)^2 + \frac{1}{2\mu} \sum_{i \in \mathcal{E}} (c_i(x))^2$$

- And I can solve it as a sequence of unconstrained problems!

---

## 9.2 MERIT FUNCTIONS AND FILTERS

# Feasible algorithms

---

- If I can afford to maintain feasibility at all steps, then I just monitor decrease in objective function.
- I accept a point if I have enough descent.
- But this works only for very particular constraints, such as linear constraints or bound constraints (and we will use it).
- Algorithms that do that are called **feasible algorithms**.

# Infeasible algorithms

---

- But, sometimes it is VERY HARD to enforce feasibility at all steps (e.g. nonlinear equality constraints).
- And I need feasibility only in the limit; so there is benefit to allow algorithms to move on the outside of the feasible set.
- But then, how do I measure progress since I have two, apparently contradictory requirements:
  - Reduce infeasibility (e.g.  $\sum_{i \in \mathcal{E}} |c_i(x)| + \sum_{i \in \mathcal{I}} \max\{-c_i(x), 0\}$  )
  - Reduce objective function.
  - It has a multiobjective optimization nature!

---

## 9.2.1 MERIT FUNCTIONS

# Merit function


---

- One idea also from multiobjective optimization: minimize a weighted combination of the 2 criteria.

$$\phi(x) = w_1 f(x) + w_2 \left[ \sum_{i \in \mathcal{E}} |c_i(x)| + \sum_{i \in \mathcal{I}} \max\{-c_i(x), 0\} \right]; \quad w_1, w_2 > 0$$

- But I can scale it so that the weight of the objective is 1.
- In that case, the weight of the infeasibility measure is called “penalty parameter”.
- I can monitor progress by ensuring that  $\phi(x)$  decreases, as in unconstrained optimization.

# Nonsmooth Penalty Merit Functions

$$\phi_1(x; \mu) = f(x) + \mu \sum_{i \in \mathcal{E}} |c_i(x)| + \mu \sum_{i \in \mathcal{I}} [c_i(x)]^-, \quad [z]^- = \max\{0, -z\}.$$


- It is called the  $\ell_1$  merit function.
- Sometimes, they can be even EXACT.

**Definition 15.1** (Exact Merit Function).

*A merit function  $\phi(x; \mu)$  is exact if there is a positive scalar  $\mu^*$  such that for any  $\mu > \mu^*$ , any local solution of the nonlinear programming problem (15.1) is a local minimizer of  $\phi(x; \mu)$ .*

We show in Theorem 17.3 that, under certain assumptions, the  $\ell_1$  merit function  $\phi_1(x; \mu)$  is exact and that the threshold value  $\mu^*$  is given by

$$\mu^* = \max\{|\lambda_i^*|, i \in \mathcal{E} \cup \mathcal{I}\},$$



# Smooth and Exact Penalty Functions

---

- Excellent convergence properties, but very expensive to compute.
- Fletcher's augmented Lagrangian:

$$\phi_F(x; \mu) = f(x) - \lambda(x)^T c(x) + \frac{1}{2}\mu \sum_{i \in \mathcal{E}} c_i(x)^2,$$

$$\lambda(x) = [A(x)A(x)^T]^{-1} A(x) \nabla f(x).$$

- It is both smooth and exact, but perhaps impractical due to the linear solve.

# Augmented Lagrangian

---

- Smooth, but inexact.

$$\phi(x) = f(x) - \sum_{i \in \mathcal{E}} \lambda_i c_i(x) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} c_i^2(x) \Rightarrow$$

- An update of the Lagrange Multiplier is needed.
- We will not use it, except with Augmented Lagrangian methods themselves.

# Line-search (Armijo) for Nonsmooth Merit Functions

---

$$\phi_1(x; \mu) = f(x) + \mu \sum_{i \in \mathcal{E}} |c_i(x)| + \mu \sum_{i \in \mathcal{I}} [c_i(x)]^-,$$

- How do we carry out the “progress search”?
- That is the line search or the sufficient reduction in trust region?
- In the unconstrained case, we had

$$f(x_k) - f(x_k + \beta^m d_k) \geq -\rho \beta^m \nabla f(x_k)^T d_k; \quad 0 < \beta < 1, 0 < \rho < 0.5$$

- But we cannot use this anymore, since the function is not differentiable.

# Directional Derivatives of Nonsmooth Merit Function

---

$$\phi_1(x; \mu) = f(x) + \mu \sum_{i \in \mathcal{E}} |c_i(x)| + \mu \sum_{i \in \mathcal{I}} [c_i(x)]^-,$$

- Nevertheless, the function has a directional derivative (follows from properties of max function). **EXPAND**

$$D(\phi(x, \mu); p) = \lim_{t \rightarrow 0, t > 0} \frac{\phi(x + tp, \mu) - \phi(x, \mu)}{t}; \quad D(\max\{f_1, f_2\}, p) = \max\{\nabla f_1 p, \nabla f_2 p\}$$

- Line Search:  $\phi(x_k, \mu) - \phi(x_k + \beta^m p_k, \mu) \geq -\rho \beta^m D(\phi(x_k, \mu), p_k);$

- Trust Region

$$\phi(x_k, \mu) - \phi(x_k + \beta^m p_k, \mu) \geq -\eta_1 (m(0) - m(p_k));$$
$$0 < \eta_1 < 0.5$$

# And .... How do I choose the penalty parameter?

---

- VERY tricky issue, highly dependent on the penalty function used.
- For the l1 function, guideline is:

$$\mu^* = \max\{|\lambda_i^*|, i \in \mathcal{E} \cup \mathcal{I}\},$$

- But almost always adaptive. Criterion: If optimality gets ahead of feasibility, make penalty parameter more stringent.
- E.g l1 function: the max of current value of multipliers plus safety factor (EXPAND)

---

## **9.2.2 FILTER APPROACHES**

# Principles of filters

---

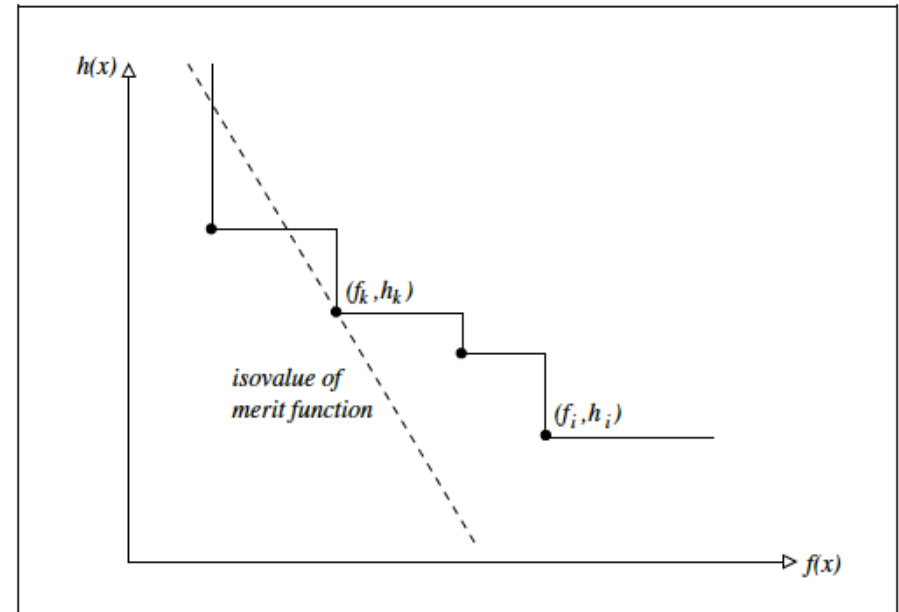
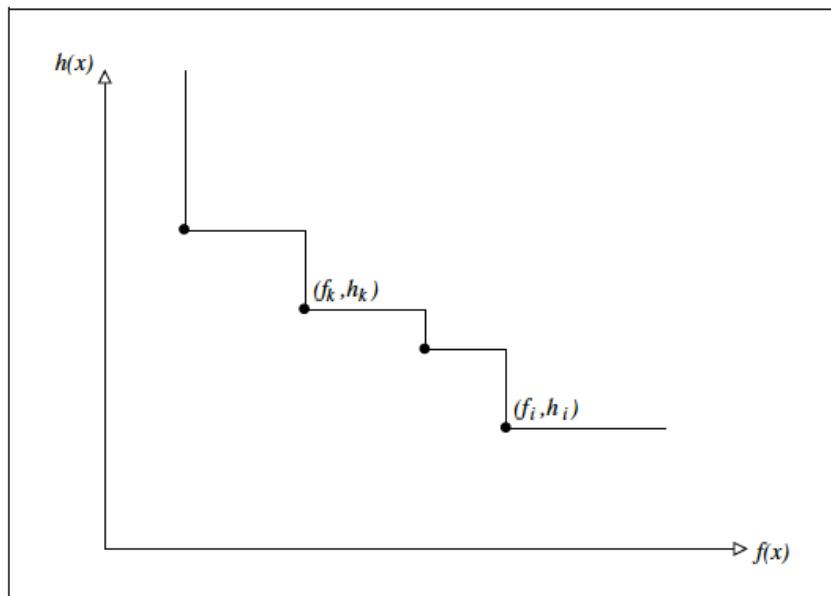
- Originates in the multiobjective optimization philosophy: objective and infeasibility

$$h(x) = \sum_{i \in \mathcal{E}} |c_i(x)| + \sum_{i \in \mathcal{I}} [c_i(x)]^-,$$

- The problem becomes:

$$\min_x f(x) \quad \text{and} \quad \min_x h(x).$$

# The Filter approach



## Definition 15.2.

- (a) A pair  $(f_k, h_k)$  is said to dominate another pair  $(f_l, h_l)$  if both  $f_k \leq f_l$  and  $h_k \leq h_l$ .
- (b) A filter is a list of pairs  $(f_l, h_l)$  such that no pair dominates any other.
- (c) An iterate  $x_k$  is said to be acceptable to the filter if  $(f_k, h_k)$  is not dominated by any pair in the filter.



## Some Refinements

---

- Like in the line search approach, I cannot accept EVERY decrease since I may never converge.
- Modification:

A trial iterate  $x^+$  is acceptable to the filter if, for all pairs  $(f_j, h_j)$  in the filter, we have that

$$f(x^+) \leq f_j - \beta h_j \quad \text{or} \quad h(x^+) \leq h_j - \beta h_j, \quad \beta \sim 10^{-5} \quad (15.33)$$

---

## **9.3 MARATOS EFFECT AND CURVILINEAR SEARCH**

# Unfortunately, the Newton step may not be compatible with penalty

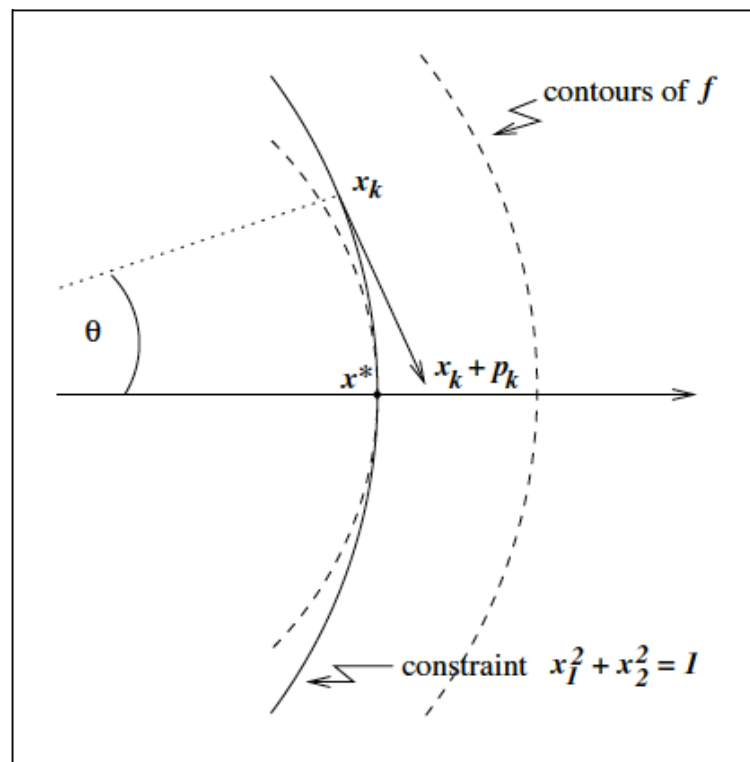
- This is called the Maratos effect.

- Problem:

$$\min f(x_1, x_2) = 2(x_1^2 + x_2^2 - 1) - x_1,$$

$$x_1^2 + x_2^2 - 1 = 0.$$

- Note: the closest point on search direction (Newton) will be rejected !
- So fast convergence does not occur



# Solutions?

---

- Use Fletcher's function that does not suffer from this problem.
- Following a step:  $A_k p_k + c(x_k) = 0$ .
- Use a correction that satisfies  $A_k \hat{p}_k + c(x_k + p_k) = 0$ .

$$\hat{p}_k = -A_k^T (A_k A_k^T)^{-1} c(x_k + p_k),$$

- Followed by the update or line search:

$$x_k + p_k + \hat{p}_k \quad x_k + \tau p_k + \tau^2 \hat{p}_k$$

- Since  $c(x_k + p_k + \hat{p}_k) = O(\|x_k - x^*\|^3)$  compared to  $c(x_k + p_k) = O(\|x_k - x^*\|^2)$  corrected Newton step is likelier to be accepted.



THE UNIVERSITY OF  
CHICAGO

# Section 10: Quadratic Programming

Reference: Chapter 16, Nocedal and  
Wright.

---

# 10.1 GRADIENT PROJECTIONS FOR QPS WITH BOUND CONSTRAINTS

# Projection

---

$$\min_x \quad q(x) = \frac{1}{2}x^T Gx + x^T c$$

- The problem: subject to  $l \leq x \leq u$ ,
- Like in the trust-region case, we look for a Cauchy point, based on a projection on the feasible set.
- $G$  does not have to be psd (essential for AugLag)
- The projection operator:

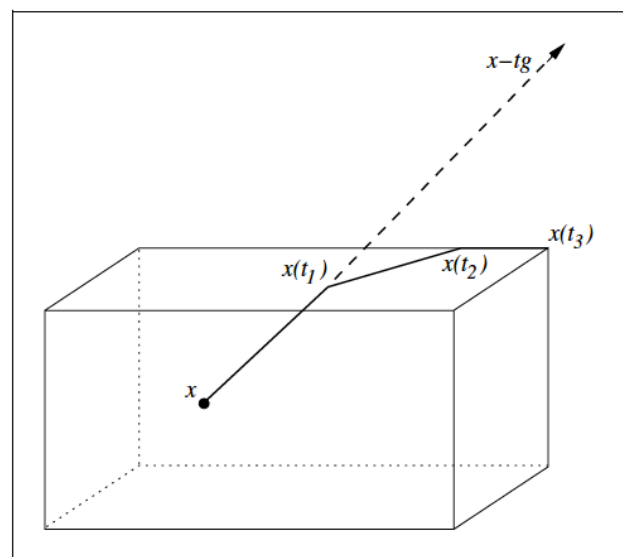
$$P(x, l, u)_i = \begin{cases} l_i & \text{if } x_i < l_i, \\ x_i & \text{if } x_i \in [l_i, u_i], \\ u_i & \text{if } x_i > u_i. \end{cases}$$

# The search path

- Create a piecewise linear path which is feasible (as opposed to the linear one in the unconstrained case) by projection of gradient.

$$x(t) = P(x - tg, l, u),$$

$$g = Gx + c;$$





# Computation of breakpoints

---

- Can be done on each component individually

$$\bar{t}_i = \begin{cases} (x_i - u_i)/g_i & \text{if } g_i < 0 \text{ and } u_i < +\infty, \\ (x_i - l_i)/g_i & \text{if } g_i > 0 \text{ and } l_i > -\infty, \\ \infty & \text{otherwise.} \end{cases}$$

- Then the search path becomes on each component:

$$x_i(t) = \begin{cases} x_i - t g_i & \text{if } t \leq \bar{t}_i, \\ x_i - \bar{t}_i g_i & \text{otherwise.} \end{cases}$$

# Line Search along piecewise linear path

---

- Reorder the breakpoints eliminating duplicates and zero values to get

$$0 < t_1 < t_2 < \dots$$

- The path:

$$x(t) = x(t_{j-1}) + (\Delta t)p^{j-1}, \quad \Delta t = t - t_{j-1} \in [0, t_j - t_{j-1}],$$

- Whose direction is:

$$p_i^{j-1} = \begin{cases} -g_i & \text{if } t_{j-1} < \bar{t}_i, \\ 0 & \text{otherwise.} \end{cases}$$

## Line Search (2)

---

- Along each piece,  $[t_{j-1}, t_j]$  find the minimum of the quadratic

$$\frac{1}{2}x^T Gx + c^T x$$

- This reduces to analyzing a one dimensional quadratic form of  $t$  on an interval.
- If the minimum is on the right end of interval, we continue.
- If not, we found the local minimum and the Cauchy point.

# Subspace Minimization

---

- Active set of Cauchy Point

$$\mathcal{A}(x^c) = \{i \mid x_i^c = l_i \text{ or } x_i^c = u_i\}.$$

- Solve subspace minimization problem

$$\begin{aligned} \min_x q(x) &= \frac{1}{2}x^T Gx + x^T c \\ \text{subject to } & x_i = x_i^c, \ i \in \mathcal{A}(x^c), \\ & l_i \leq x_i \leq u_i, \ i \notin \mathcal{A}(x^c). \end{aligned}$$

- No need to solve exactly. For example truncated CG with termination if one inactive variable reaches bound.

# Gradient Projection for QP

---

**Algorithm 16.5** (Gradient Projection Method for QP).

Compute a feasible starting point  $x^0$ ;

for  $k = 0, 1, 2, \dots$

if  $x^k$  satisfies the KKT conditions for (16.68)

stop with solution  $x^* = x^k$ ;

Set  $x = x^k$  and find the Cauchy point  $x^c$ ;

Find an approximate solution  $x^+$  of (16.74) such that  $q(x^+) \leq q(x^c)$   
and  $x^+$  is feasible;

$x^{k+1} \leftarrow x^+$ ;

end (for)



Or, equivalently, if projection does not advance from 0.